



LAWRENCE
LIVERMORE
NATIONAL
LABORATORY

Global metagenomic survey reveals taxonomic "blind spots" in geothermal springs

E. A. Elie-Fadrosh, D. Paez-Espino, J. Jarett, P. F. Dunfield, B. P. Hedlund, A. E. Dekas, S. E. Grasby, A. L. Brady, H. Dong, B. R. Briggs, W. Li, D. Goudeau, R. Malmstrom, A. Pati, J. Pett-Ridge, E. M. Rubin, T. Woyke, N. C. Kyrpides, N. N. Ivanova

August 18, 2015

Nature Communications

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

1 **LLNL-JRNL-676401**

2 **Global metagenomic survey reveals a new bacterial candidate phylum in**
3 **geothermal springs**

4
5 Emiley A. Eloie-Fadros¹, David Paez-Espino¹, Jessica Jarett¹, Peter F. Dunfield², Brian
6 P. Hedlund³, Anne E. Dekas⁴, Stephen E. Grasby⁵, Allyson L. Brady⁶, Hailiang Dong⁷,
7 Brandon R. Briggs⁸, Wen-Jun Li⁹, Danielle Goudeau¹, Rex Malmstrom¹, Amrita Pati¹,
8 Jennifer Pett-Ridge⁴, Edward M. Rubin^{1,10}, Tanja Woyke¹, Nikos C. Kyrpides¹, Natalia
9 N. Ivanova^{1*}

10
11 **Affiliations:**

12 ¹Joint Genome Institute, Walnut Creek, CA 94598, USA.

13 ²University of Calgary, Calgary, AB T2N 1N4, Canada.

14 ³University of Nevada, Las Vegas, Las Vegas, NV 89154, USA.

15 ⁴Lawrence Livermore National Laboratory, Livermore, CA 94550, USA.

16 ⁵Geological Survey of Canada, Calgary, AB, Canada.

17 ⁶McMaster University, Hamilton, ON L8S 4L8, Canada.

18 ⁷Miami University, Oxford, OH, 45056, USA.

19 ⁸University of Alaska-Anchorage, Anchorage, AK, USA.

20 ⁹Sun Yat-Sen University, Guangzhou, 510275, China.

21 ¹⁰Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA.

22
23 *Correspondence to: nnivanova@lbl.gov.

Abstract

Analysis of the increasing wealth of metagenomic data collected from diverse environments can lead to the discovery of novel branches on the tree of life. Here we analyse 5.2 Tb of metagenomic data collected globally to discover a novel bacterial phylum ('*Candidatus* Kryptonia') found exclusively in high-temperature pH-neutral geothermal springs. This lineage had remained hidden as a taxonomic "blind spot" due to mismatches in the primers commonly used for ribosomal gene surveys. Genome reconstruction from metagenomic data combined with single-cell genomics results in several high-quality genomes representing four genera from the new phylum. Metabolic reconstruction indicates a heterotrophic lifestyle with conspicuous nutritional deficiencies, suggesting the need for metabolic complementarity with other microbes. Co-occurrence patterns identifies a number of putative partners, including an uncultured *Armatimonadetes* lineage. The discovery of Kryptonia within previously studied geothermal springs underscores the importance of globally sampled metagenomic data in detection of microbial novelty, and highlights the extraordinary diversity of microbial life still awaiting discovery.

47

48 Molecular environmental surveys have provided a sizeable snapshot of microbial
49 phylogenetic diversity. Sequencing of small subunit ribosomal RNA (SSU rRNA) genes
50 directly from the environment has expanded the known microbial tree of life from
51 Woese's original twelve phyla to more than 70 bacterial phyla^{1,2}. Advances in
52 cultivation-independent methods for examining uncultured microbes, including single-
53 cell genomics and deep sequencing of environmental samples, have begun yielding
54 complete or near-complete genomes from many novel lineages³⁻¹⁰. These approaches
55 have already led to the recovery of genomic information from a wealth of candidate
56 lineages (phylogenetic lineages for which a cultured representative is not available),
57 notably the Lokiarchaeota¹¹, Pacearchaeota and Woesearchaeota¹⁰, and members of
58 the Candidate Phyla Radiation³. These lineages, previously recognized only through
59 SSU rRNA data and residing in poorly sampled habitats, are providing a more complete
60 topology of the tree of life.

61 More recently, it has been suggested that a wealth of novel bacterial and
62 archaeal clades exist that are systematically under-represented (the 'rare biosphere') or
63 missed altogether in classical surveys, leaving significant taxonomic "blind spots"¹².
64 Compared to many of the proposed candidate phyla for which SSU rRNA gene
65 information exists, these taxonomic "blind spots" are uncharted lineages with potentially
66 important ecological and evolutionary implications. Further, these lineages may be
67 highly abundant and hold important metabolic or functional roles within the community,
68 yet have been overlooked thus far in ecological surveys. Metagenome sequencing is
69 uniquely suited for uncovering taxonomic "blind spots" since it does not suffer from

biases introduced during PCR amplification, and has limitations only with insufficient resolution of minor populations within a community. However, an exploration of the complete compendium of available metagenomic sequences for the presence of taxonomic "blind spots" has yet to be performed¹³. Here, we report the results of large-scale mining of metagenomic data and single cell genomics, which led to the discovery of a new bacterial phylum in geographically distinct geothermal springs.

Results

Identification of a novel bacterial candidate phylum

To cast a global net for the discovery of novel microbial lineages in the absence of biases introduced via PCR amplicon-based surveys, we collected long assembled contigs (≥ 100 kbp) from a comprehensive collection of 4,290 metagenomic datasets available through the Integrated Microbial Genomes with Microbiome Samples (IMG/M), a database containing a total of more than 5 Tb of sequence data¹⁴. From these data, 31,955 assembled contigs were identified and 744 contigs were further selected that contained SSU rRNA gene fragments greater than 100 bp (Fig. 1A). The SSU rRNA gene sequences were then aligned and phylogenetically placed on a reference tree consisting of high-quality SSU rRNA sequences from Bacteria and Archaea^{15,16}. Exploration of the constructed SSU rRNA tree for novel phylogenetic branches led to the identification of a distinct lineage consisting of a full-length SSU rRNA sequence. A subsequent search against all assembled metagenomic data identified three additional full-length SSU rRNA sequences. The four SSU rRNA gene sequences were from four geographically distant, high-temperature, pH-neutral, geothermal springs in North

America and Asia (Fig. 1). These sequences shared an average 97.4% identity (\pm 1.97% s.d.), and showed a maximum identity of only 83% to SSU rRNA genes (such as the one in GenBank ID: AP011715) in NCBI's Non-Redundant (NR) database. In line with the notion of taxonomic "blind spots"¹², a comparison of 'universal' SSU rRNA primer sets typically used for full-length and hypervariable region amplification with the four novel sequences indicated numerous mismatches, explaining why members of this lineage likely eluded detection in previous microbial diversity surveys (Supplementary Fig. 1; Supplementary Table 1).

Phylogenetic analysis of the four SSU rRNA genes placed the newly discovered lineage into a monophyletic branch within the *Fibrobacteres-Chlorobi-Bacteroidetes* (FCB) superphylum^{9,17} (Supplementary Fig. 2). Based on suggested thresholds for SSU rRNA sequence identity to distinguish new phyla^{2,18}, we propose that this lineage represents a new bacterial candidate phylum (Supplementary Table 2).

Comparative genomics and cell morphology of novel FCB lineage

Reassembly of the metagenomic data combined with tetranucleotide-based binning methods using the initial contigs containing the SSU rRNA genes yielded near-complete recovery of four distinct genomes, each from one of the four spring samples (Supplementary Fig. 3; Supplementary Table 3). Phylogenetic analysis of conserved marker genes supported its placement as a sister phylum to the *Ignavibacteria* with 100% bootstrap support (Fig. 2A; Supplementary Fig. 4). Three of the genomes reconstructed from metagenomes (GFMs) from Dewar Creek Spring, Canada¹⁹, Great Boiling Spring, Nevada^{20,21}, and Gongxiaoshe pool, Yunnan Province, China²² had an

average 95.8% estimated coverage, while the genome from Jinze pool, Yunnan Province, China²² had a lower estimated coverage of 68% (Supplementary Table 4). The high genomic sequence coverage across the four metagenomes (average 31.2x coverage; Supplementary Table 3) suggested that this novel lineage might exist at sufficient cell abundance to be captured by single cell technology. We therefore employed high-throughput single-cell isolation, whole-genome amplification and SSU rRNA screening of single amplified genomes (SAGs) in search for the novel lineage (Fig. 1). We successfully recovered a total of 18 SAGs from three of the four samples, corresponding to the novel phylum-level clade with an estimated average genome completeness of 67.2% (± 20.1 s.d.) (Supplementary Table 3). We designate this new candidate phylum '*Candidatus* Kryptonia,' from the Greek word '*krupton*' meaning hidden or secret since it has hitherto eluded detection due to SSU rRNA primer biases (Supplementary Table 4).

The average nucleotide identity (ANI) based metric, Microbial Species Identifier (MiSI), was used to compare the four '*Ca. Kryptonia*' genomes reconstructed from metagenomes (GFM) and the 18 SAGs²³. This analysis revealed that almost all of the genotypes extracted from the same sample belonged to a single species (Supplementary Table 5). For example, the GFM reconstructed from Dewar Creek ('*Ca. Kryptonium thompsoni*' JGI-4) and the thirteen SAGs ('*Ca. Kryptonium thompsoni*' JGI-5 – JGI-17) collected from the same site shared an ANI of 99.67% (± 0.15 s.d.) and represent a single coherent species²³. A single exception to the above observations was the recovery of a divergent '*Ca. Kryptonia*' SAG ('*Ca. Chrysopegis kryptomonas*' JGI-23) from the Jinze pool, Yunnan Province, China representing a population distinct from

the other two SAGs recovered from this site ('*Ca. Kryptobacter tengchongensis*' JGI-24 and JGI-25) (Supplementary Table 5). Across the four geothermal springs, the GFMs and SAGs collectively share average ANIs of only 78.86% (± 1.42 s.d.), suggesting that they represent different genera of 'Ca. Kryptonia'. Further support for genus-level designations is evident from nuanced functional and metabolic differences across the genomes, as described below.

In addition to recovering single cells of *Ca. Kryptonia* for genome amplification, we designed a SSU rRNA-targeted fluorescence *in situ* hybridization (FISH) probe to visualize cell morphology (Fig. 2B). The targeted *Ca. Kryptonia* cells appeared filamentous, and exhibited morphological heterogeneity ranging from short to elongated filaments. These findings are consistent with numerous reports describing filamentous thermophilic bacteria, most notably cultivated members of the sister phylum *Ignavibacteria* that range in length from 1 μm to greater than 15 μm ^{24,25}.

Unique CRISPR-Cas fusion and limited biogeographic distribution

CRISPR (clustered regularly interspaced short palindromic repeats) elements and *cas* (CRISPR-associated) genes across the *Ca. Kryptonia* genomes were recovered, and are suggestive of defense against viral attack. A novel fusion between two different CRISPR-Cas types (type I and III; subtypes I-B and III-A) was identified in all genomes. This unusual fusion contained the full gene set for components responsible for the multistep CRISPR processes for spacer acquisition, CRISPR locus transcription and maturation, and final nucleic acid interference^{26,27} (Supplementary Fig. 5). This observation represents the first report of a type I-B/type III-A CRISPR-Cas

fusion and expands the known genetic diversity of CRISPR-Cas loci. Based on reconstruction of repeat-spacer arrays, the *Ca. Kryptonium thompsoni* genomes appear to represent a clonal CRISPR population without active spacer acquisition, while the *Ca. Kryptobacter tengchongensis* genomes are considerably dynamic in terms of a mosaic spacer collection (Supplementary Text; Supplementary Tables 6 and 7). These findings suggest that the CRISPR-Cas encoded by *Ca. Kryptobacter tengchongensis* is highly active, while the *Ca. Kryptonium thompsoni* genomes are not actively acquiring spacers through the CRISPR-Cas system.

To verify the limited biogeographic distribution of *Ca. Kryptonia*, we systematically surveyed the collection of 640 Gb of assembled metagenomic data from 4,290 environmental samples (including 169 samples from geothermal springs and hydrothermal vents) for the presence of a genomic signature beyond our initial search using SSU rRNA fragments from 100 kbp contigs (Fig. 3; Supplementary Table 8). Further, we searched against all available SSU rRNA data from the SILVA database¹⁶ for additional *Ca. Kryptonia* phylotypes and did not recover a highly similar match. Using this expanded search, we found evidence for *Ca. Kryptonia* in a total of twenty metagenomes, which included only three additional geographic sites compared to our initial SSU rRNA survey (Supplementary Table 8). The environments where this phylum was found were similar to the settings where we first discovered the genomic presence of *Ca. Kryptonia*: all were high-temperature ($\geq 70^{\circ}\text{C}$), pH-neutral (6.4 – 8.0) settings. In sum, the limited range of *Ca. Kryptonia* is reflected in the observation that genomic signatures were found in nine unique geographical locations from a total of twenty-three

pH neutral hot springs currently sampled by metagenomics, and absent from the 1,614 unique locations represented by 4,290 metagenomic samples.

Additional metagenomic searches specific for all CRISPR repeat-spacer arrays collected from the *Ca. Kryptonia* genomes resulted in a similar pattern of limited biogeographic distribution (Fig. 3; Supplementary Table 7). We identified shared spacers across *Ca. Kryptonia* populations in geographically distinct geothermal springs. For example, shared spacers were identified between the *Ca. Kryptobacter tengchongensis* JGI-2 and JGI-3 genomes despite sampling from separate geothermal pools in China. Further, shared spacers were identified across exceptionally wide geographic distances including Canada and Nevada (*Ca. Kryptonium thompsoni* JGI-4 and the Great Boiling Springs metagenome), and China and Nevada (*Ca. Kryptobacter tengchongensis* JGI-2 and the Great Boiling Springs metagenome) (Fig. 3). Remarkably, we also found spacer matches to a set of metagenomic contigs that we assigned as viral due to their linkage to known viral genes, from these same samples and metagenome samples collected from Yellowstone National Park²⁸ (Fig. 3; Supplementary Text; Supplementary Fig. 5; Supplementary Table 9). These genomic recruitment and spacer signature data suggest that *Ca. Kryptonia* is present in additional geothermal spring sites and that viruses which appear to infect *Ca. Kryptonia* circulate across wide geographic space as revealed from the conserved infection vestiges.

Metabolic potential of ‘Candidatus Kryptonia’

The availability of multiple nearly complete *Ca. Kryptonia* genomes from both GFMs and SAGs enabled metabolic and putative functional predictions for this novel candidate phylum, as well as insights into some of the unique properties and notable absence of function for the individual genera. Approximately 50% of the predicted composite proteome for the *Ca. Kryptonia* genomes showed similarity to a diverse array of FCB superphylum members, with 11.3% and 1.96% best matches to thermophilic members of the phylum *Ignavibacteria* and *Caldithrix abyssi*, respectively (Supplementary Fig. 6). The conserved Por secretion system C-terminal sorting domain (TIGR04183), found exclusively in members of the FCB superphylum⁹, was recovered in all GFMs and SAGs, and altogether totaled 811 predicted proteins across the *Ca. Kryptonia* genomes. Reverse gyrase, the presumptive gene indicator for the extreme thermophilic and hyperthermophilic lifestyle in bacteria and archaea²⁹, was found in all *Ca. Kryptonia* genomes, which suggests that most, if not all members, of this lineage are extreme thermophiles or hyperthermophiles. Further, we found evidence for horizontal gene transfer of the reverse gyrase from the crenarchaeal order *Thermoproteales* (Supplementary Text; Supplementary Fig. 7) and hypothesize that *Ca. Kryptonia*'s thermophilic traits might have been acquired via lateral gene transfer rather than ancestral inheritance.

Ca. Kryptonia is a motile heterotroph with a complete tricarboxylic acid cycle and key metabolic enzymes for Embden-Meyerhof glycolysis and the pentose phosphate pathway. We found evidence for a complex oxidative phosphorylation pathway, which points towards aerobic respiration (Fig. 4; Supplementary Table 10). An elaborate and unique respiratory pathway for the redox transformation of iron is encoded in the *Ca.*

Kryptonite genomes with similar, yet non-homologous components to the well-characterized Mtr-like respiratory pathway³⁰ (Supplementary Fig. 8). Altogether, *Ca. Kryptonite* has the machinery to carry out ferric iron respiration under thermophilic conditions and likely vies with archaeal community members to impact metal biogeochemistry in these geothermal springs.

Ca. Kryptonite hosts the genomic potential for aromatic hydrocarbon degradation via oxidation to catechol, and subsequent catechol meta-cleavage (Fig. 4). Further, the *Ca. Kryptonite thompsoni* genomes encode a putative gene complement for the anaerobic degradation of aromatic amino acids or similar compounds, notably represented by a phenylacetyl-CoA oxidoreductase homologous to the hyperthermophilic archaeon *Ferroglobus placidus*³¹. This feature appears to be the first example of an extremely thermophilic or hyperthermophilic bacterium with the presumptive capacity to completely mineralize aromatic compounds, and holds biotechnological potential as well as implications for carbon cycling within geothermal springs³².

Unexpected metabolic deficiencies identified in Ca. Kryptonite

An unexpected observation was that all *Ca. Kryptonite* genomes had conspicuous nutritional deficiencies, displaying gene loss for many biosynthetic pathways, including thiamine, biotin and amino acids, such as the evolutionarily conserved histidine biosynthesis³³ (Fig. 4; Supplementary Table 10). While obligately host-dependent microbes and some free-living organisms with reduced genomes are known to omit a suite of anabolic pathways^{34,35}, the *Ca. Kryptonite* genomes do not appear to have

signatures of either lifestyle. An analysis of 759 high-quality FCB superphylum genomes indicate the near-complete *Ca. Kryptonia* genomes are distinct from free-living microbes in terms of amino acid pathway coverage and genome size, yet are not highly reduced compared to obligate symbionts (Supplementary Fig. 9). These findings suggest that *Ca. Kryptonia* has potentially evolved functional dependency on other microbes in order to acquire necessary metabolic requirements.

To explore the existence of possible microbial partners, we performed a co-occurrence analysis of SSU rRNA sequences retrieved through their targeted assembly from an expanded set of 22 geothermal springs metagenomes (Supplementary Text; Supplementary Table 11). An analysis of co-occurrence patterns for clusters of taxonomically coherent groups (clustered at 90% sequence identity) revealed a subset of taxonomically clustered groups (phylotypes) highly correlated with the abundance of *Ca. Kryptonia* (Supplementary Table 12). These clusters included an *Armatimonadetes* lineage, which had the highest correlation value, three separate lineages of *Chloroflexi*, and *Thermus* spp. (Fig. 5). For the twelve metagenomes in which *Ca. Kryptonia*'s SSU rRNA was reconstructed, the *Armatimonadetes* lineage was found to co-occur in seven of those metagenomes at similar sequence coverage to the *Ca. Kryptonia* genomes, and was conspicuously absent across all other metagenomes surveyed. To explore the potential of the *Armatimonadetes* lineage to complement the metabolic deficiencies identified in *Ca. Kryptonia*, we reconstructed three nearly complete genomes of *Armatimonadetes* (Fig. 2; Supplementary Tables 3 and 13) to infer metabolic potential and signatures of possible metabolic exchange and interaction. Analysis of the reconstructed genomes identified metabolic features complementary to those of *Ca.*

Kryptonite, such as histidine, cysteine and methionine, proline, aspartic acid, and thiamine biosynthesis, and degradation of pentoses (Fig. 5B; Supplementary Text; Supplementary Table 13). Furthermore, in the reconstructed *Armatimonadetes* genomes we also identified a CsgG family protein, which forms transmembrane channels for secretion of “functional amyloids,” a class of bacterial proteins capable of assembling highly stable fibers through a nucleation-precipitation mechanism³⁶. “Functional amyloids” play major roles in adhesion to surfaces and biofilm formation in diverse bacteria including *Escherichia coli*, *Caulobacter crescentus* and *Bacillus subtilis*³⁷. Further, the CsgG-like transporter was located in a six-gene conserved cluster containing a predicted subtilase-family peptidase and a putative secreted protein with four copies of a “carboxypeptidase regulatory-like domain” (Pfam13620) (Supplementary Fig. 10). This domain is a member of the transthyretin clan and has been found to form amyloid in physiological conditions³⁸. We hypothesize that this cluster in the *Armatimonadetes* genomes encodes for synthesis, secretion and assembly of “functional amyloid,” in which other members of the community may be embedded. On the other hand, the *Ca. Kryptonite* genomes encode many proteases and peptidases, which may be responsible for remodeling and digestion of this extracellular matrix.

Other co-occurring lineages with *Ca. Kryptonite* include the *Thermus* spp. cluster (Supplementary Table 12). Interestingly, *Ca. Kryptonite* might complement an incomplete denitrification pathway in *Thermus* spp., which may be responsible for high rates of nitrous oxide production^{39,40}. *Thermus* spp. have been experimentally characterized to reduce nitrate to nitrous oxide but lack the capacity to subsequently

produce dinitrogen^{39,40}. *Kryptonia* encodes a nitrous oxide reductase (EC 1.7.2.4) but lacks other components of the denitrification pathway (Supplementary Text; Supplementary Table 14). Taken together, we hypothesize that *Ca. Kryptonia* may participate in a partnership with other organisms, such as the *Armatimonadetes*, or might interact with a broader consortium of microbes within the geothermal spring environment.

Discussion

A comprehensive survey of a global set of assembled metagenomic data for novel microbial lineages has resulted in the discovery of a new bacterial candidate phylum in geothermal springs. The high-quality draft genome assemblies enabled by complementary approaches from metagenomic data and single-cell genomics data for *Ca. Kryptonia* facilitated delineation of the host-virus interaction across geographically distant sites. Further, we observed a novel fusion between two different CRISPR-Cas types, representing the first report of a type I-B/type III-A CRISPR-Cas fusion and expanded the known genetic diversity of CRISPR-Cas loci.

The metabolic capacity for *Ca. Kryptonia* provides evidence for a unique heterotrophic lifestyle with the putative capacity for iron respiration within a consistent ecological niche in geothermal springs. An unexpected observation was that all *Ca. Kryptonia* genomes had conspicuous nutritional deficiencies, which led to the hypothesis of a microbial partnership or interaction with a broader consortium of microbes. Subsequent genome reconstruction of genomes from a co-occurring *Armatimonadetes* lineage indicated potential complementarity for those metabolic

features presumably absent in *Ca. Kryptonia*. It is well recognized that certain marine microbes, such as SAR11 (ref.⁴¹) and SAR86 (ref.⁴²), lack a variety of anabolic pathways and likely rely on other microbial community members to supplement their requirements. Within geothermal springs, the growth of chlorophototroph *Candidatus Chloracidobacterium thermophilum* in the laboratory was shown to depend upon two heterotrophs, *Anoxybacillus* and *Meiothermus* spp., due to lack of biosynthetic pathways for branched-chain amino acids, lysine and cobalamin⁴³. Our study suggests that dependency on other organisms within the geothermal spring community might be a more common occurrence than previously appreciated, perhaps contributing to challenges in obtaining many of these lineages as isolated monocultures. Future efforts to delineate this hypothesized interaction, particularly utilizing microscopy methods to visualize these uncultivated cells *in situ*, will further contribute to our understanding of *Ca. Kryptonia* and its role within the environment.

Geothermal springs have been heavily surveyed as a rich source of novel microbial branches on the tree of life^{18,44}, yet our results indicate that additional phylogenetic novelty has yet to be captured from these environments. The discovery of a new candidate phylum emphasize that extraordinary microbial novelty is likely still awaiting discovery using the vast metagenomic data assembled from locations sampled globally.

Materials and methods

Metagenomes

All publicly available metagenome datasets from IMG/M were used in the study¹⁴. The metagenomes can be accessed at <http://img.jgi.doe.gov> and associated metadata can be found in the GOLD database at <http://genomesonline.org>.

Metagenomic binning

Tetranucleotide-based binning methods were implemented as previously described to recover near-complete genomes from metagenomes⁴⁵. Both single metagenomes and combined metagenome assemblies were used to recruit additional contigs that harbored the same tetranucleotide signature, and the raw reads were subsequently re-assembled using SPAdes version 3.1.0 (ref.⁴⁶).

SAG generation

Sediment samples were collected from Dewar Creek hot spring (49.9543667°, -116.5155000°) near the source of the hot spring on September 28, 2012, from the Jinze pool (25.44138°, 98.46004°) on August 12, 2012, and from the Gongxiaoshe pool (25.44012°, 98.44081°) on August 9, 2011. Samples were mixed with glycerol-TE buffer for cryopreservation and stored at -80°C. Single cells were isolated using FACS, lysed, and subjected to whole-genome amplification (WGA) as previously described⁹ with the following modifications: the alkaline lysis was preceded by a 20 min digest with lysozyme (Epicentre) at 30°C; WGA was performed with a REPLI-g Single Cell Kit

(Qiagen) with a scaled-down reaction volume of 2 µl; and the amplification reaction was incubated for 6 hr at 30°C. WGA reactions were diluted 10-fold, then aliquots were further diluted 200-fold for PCR screening targeting the V6-V8 regions of the SSU rRNA using a QuantiNova SYBR Green PCR kit (Qiagen) for 45 cycles of amplification⁹. PCR products were purified and sequenced, and SAGs matching *Kryptonia* SSU rRNA sequences were selected for shotgun sequencing.

SAG sequencing, assembly and QC

Draft genomes for the eighteen SAGs were generated at the DOE Joint Genome Institute (JGI) using the Illumina MiSeq technology according to standard protocols (<http://www.jgi.doe.gov/>). Assembly was performed using SPAdes version 3.1.0 (ref.⁴⁶) using the --sc flag to denote MDA-derived data to account for uneven coverage of the single-cell genomes. Quality control and contaminant removal from the resultant assemblies was achieved using a two-step process. First, all assembled reads were used as input for a newly developed single-cell decontamination method (ProDeGe)⁴⁷, which uses both taxonomic and k-mer based decisions to flag putative non-target contigs. Since the taxonomic information was limited to phylum-level designations, we further supplemented this procedure with direct mapping to the genomes reconstructed from metagenomic data. For mapping, a combination of blast and blat were implemented to validate correct recruitment of the assembled SAG contigs to *Ca. Kryptonia*-specific GFM scaffolds. This method was important for retaining CRISPR/Cas genetic regions since ProDeGe had the tendency to flag these contigs based on divergent k-mer frequencies. Gene annotation was performed within the Integrated Microbial Genomes (IMG) platform developed by the DOE Joint Genome Institute¹⁴.

SSU rRNA phylogeny

Full-length SSU rRNA gene sequences from *Ca. Kryptonia* were aligned using the SINA aligner¹⁵ to a comprehensive database of references (SILVA-NR version 119)¹⁶. A total of 187 full-length bacterial and archaeal reference sequences were selected from the SILVA database, and 1,354 distinct alignment patterns were used, and filtered using the *E. coli* positional mask. A maximum likelihood tree was calculated from the masked alignments with 100 bootstrap resamplings using the Generalized Time-Reversible model with G+I options in RAxML version 7.6.3 (raxmlHPC-PTHREADS-SSE3 -f a -x 12345 -p 12345 -# 100 -T 5 -m GTRGAMMAI)⁴⁸. To resolve placement within the FCB superphylum, a subset of 77 FCB superphylum members and 37 archaeal reference sequences were selected and phylogenies constructed using two separate algorithms: the GTR+G+I model (used with maximum likelihood, RAxML⁴⁸) and Bayesian inference (MrBayes⁴⁹). Node stability was evaluated using a rapid bootstrapping analysis (RAxML, 100 runs) and posterior probabilities (MrBayes, 2.4 million generations, burnin of 25%). Alignments and phylogenetic trees are available at <http://portal.nersc.gov/dna/microbial/prokpubs/Kryptonia>.

Catalyzed Reporter Deposition Fluorescence *In Situ* Hybridization (CARD-FISH)

An oligonucleotide probe specific for *Ca. Kryptonia* (Kryp56; 5'-CCGTGTCCCTGACTTGCA-3') was designed in ARB (version 6.0.2)⁵⁰. The probe is a perfect match to 19 of the 22 *Ca. Kryptonia* SSU rRNA gene sequences recovered in this study, and contains two or more mismatches to all SSU rRNA gene sequences in the SILVA-NR database (version 123)¹⁶. The probe sequence was synthesized by Biomers.net (Ulm, Germany) with horseradish peroxidase (HRP) conjugated to the 5'

end. Cells from Dewar Creek sediment were separated from particulates by brief vortexing followed by centrifugation. Suspended cells were preserved with dimethyl sulfoxide (DMSO) and stored at -80°C. The cells were permeabilized with lysozyme and CARD-FISH was performed based on the protocol of Pernthaler and colleagues⁵¹. Hybridization was carried out at 46°C with 20% formamide, and the amplification was performed with tyramides conjugated to Alexa 488 (Life Technologies, #T20948). The optimal formamide concentration and specificity was predicted using mathFISH⁵² and the DECIPHER ProbeMelt tool⁵³ (Supplementary Table 15), and confirmed empirically by performing CARD-FISH on the Dewar Creek cells over a gradient of formamide concentrations (10 – 35%). Samples were counterstained with 4',6-diamidino-2-phenylindole (DAPI) in VECTASHEILD Antifade Mounting Media (Vector Laboratories, #H-1200). Cells were visualized and imaged using a Leica DM6000B microscope using a HCX PL APO 100X oil immersion objective.

Conserved single-copy and housekeeping gene phylogenetic inference

A set of 56 universally conserved single copy proteins in the Bacteria and Archaea was used for phylogenetic inference (Supplementary Table 16). Marker genes were detected and aligned with hmmsearch and hmalign included in HMMER3 (ref.⁵⁴) using HMM profiles obtained from phylosift (<http://phylosift.wordpress.com/>)⁵⁵. Alignments were concatenated and filtered⁵⁶. Housekeeping genes were aligned using MAFFT with mafft-linsi option⁵⁷. Best substitution model was selected using proptest⁵⁸. Phylogeny was inferred using maximum likelihood methods with RAxML (version 7.6.3)⁴⁸. Tree topologies were tested for robustness using 100 bootstrap replicates with the LG+I+G model (raxmlHPC-PTHREADS-SSE3 -f a -x 12345 -p 12345 -# 100 -m

PROTGAMMALG -T 5). Trees were visualized using Dendroscope⁵⁹. The concatenated protein alignment and phylogenetic tree is available at <http://portal.nersc.gov/dna/microbial/prokpubs/Kryptonion>.

Phylogenetic distribution of predicted proteins

The taxonomic distribution of all proteins across the genomes reconstructed from metagenomic data along with the *Ca. Kryptonion* SAGs was compiled based on best matches to a comprehensive protein database of high-quality non-redundant bacterial and archaeal isolate genomes¹⁴. This search was performed using usearch (version 7.0)⁶⁰, where a protein match was considered for proteins with $\geq 30\%$ sequence identity across $\geq 50\%$ of the query alignment length. Phylogenetic affiliation at the phylum level was assigned for top matches, while proteins lacking a match according to the above criteria were noted as 'no match.'

Biogeographic survey of *Ca. Kryptonion* against available metagenomic data

All genomic data for *Ca. Kryptonion* was searched against the assembled metagenomic data from 4,290 environmental samples using blat with the -fastMap option⁶¹. Significant matches for non-ribosomal genomic regions were considered for sequences ≥ 250 bp in length and with $\geq 75\%$ identity threshold. For metagenomic contigs mapping to the ribosomal operon, a 97% identity threshold was used to capture only high-quality matches to *Ca. Kryptonion*. Visualization of metagenomic matches globally was performed using the R package 'maps'⁶². All genomic matches can be found in Supplementary Table 8.

CRISPR-Cas locus type determination

We used 99 CRISPR-associated (*cas*) gene sequence alignments and hidden Markov models from the TIGRFAM database (originally built by Haft *et al.*⁶³ and later expanded by Zhang *et al.*⁶⁴) to precisely find and identify Cas family members within the scaffolds of the *Ca. Kryptonia* genomes. We recovered and classified the corresponding CRISPR type for complete and partial CRISPR-Cas loci in all genomes following the unified CRISPR classification from 2011 (ref.⁶⁵).

CRISPR repeat-spacer arrays analysis

The CRISPR Recognition Tool (CRT)⁶⁶ was used to detect CRISPR repeat-spacer regions across all *Ca. Kryptonia* assembled scaffolds using specific parameters (according to the JGI's annotation pipeline). In the case of *Ca. Thermokryptus mobilis* GFM JGI-1, we were unable to detect spacers, and therefore we additionally used the CRISPR assembler algorithm (Crass)⁶⁷ on the raw reads. Spacers were manually curated to cull false positives from the dataset that clearly did not represent authentic spacer regions (in sum, 38 false positives). Potentially active repeat-spacer arrays were inferred based on direct association with a *cas* gene locus. We also considered the isolated repeat-spacers arrays when they shared the same repeat sequence with associated *cas* genes. CRISPRmap^{68,69} was used to further characterize identified repeat regions. From a total of 1,031 trusted spacers, we next clustered these into 795 groups based on identity $\geq 90\%$ over the whole spacer length. Spacer groups were BLAST queried against distinct databases including *Ca. Kryptonia* genomes, reference public plasmid and viral datasets (from NCBI), and across the broad available metagenomic space (IMG/M).

Extraction, assembly and clustering of SSU rRNA genes from shotgun

metagenomic data

Raw reads aligning to 16S and 18S rRNAs were collected for 22 metagenomes (Supplementary Table 11) from geothermal environments using hmalign⁵⁴ against hmm models representing bacterial, archaeal and eukaryotic sequences⁷⁰ and also by BBmap with default settings⁷¹ against sequences from the SILVA database (version 119)¹⁶ dereplicated at 95% identity using UCLUST⁶⁰. Collected paired-end Illumina reads were merged using BBmerge⁷¹ and assembled using Newbler (v. 2.8)⁷² with -ml 60 -mi 99 -rip options. Resulting contigs and scaffolds were screened using cmalign from Infernal 1.1 package⁷³ and Rfam 16S and 18S rRNA models (RF00177.cm, RF01959.cm and RF01960.cm)⁷⁴. 16S and 18S rRNA sequences longer than 300 nt were retained and trimmed using cmalign against the best-matching model with '--matchonly' option to remove introns. Reference sequences from the SILVA database were trimmed using cmalign with a domain-specific model and '--matchonly' option, and clustered together with 16S sequences extracted from shotgun metagenome data using UCLUST and percent identity cutoffs of 94%, 92% and 90%. Clusters including sequences from at least two metagenome samples were retained and their abundances in metagenome samples were computed by multiplying the length of SSU rRNA sequence by the average coverage. Taxonomy was assigned to the clusters as last common ancestor (LCA) of SILVA reference sequences included in the cluster, or as LCA of SILVA sequences in the larger cluster obtained by co-clustering SILVA and metagenome sequences at 83% identity. Spearman's rank-order correlation of cluster abundances was used to estimate co-occurrence of the clusters in metagenome data.

503

504 Accession Codes

505 Genome sequence data have been deposited as Whole Genome Shotgun projects at
506 DDBJ/EMBL/GenBank with the accession codes LDNR000000000, LDYP000000000-
507 LDZE000000000, and LDXU000000000-LDYB000000000. The versions described in this
508 paper are versions 010000000. The annotated genomes and single-cell assemblies can
509 be accessed via IMG (<http://img.jgi.doe.gov>).

510

511 References

- 512 1. Pace, N. R. Mapping the tree of life: Progress and prospects. *Microbiol. Mol. Biol.*
513 *Rev.* **73**, 565-576 (2009).
- 514 2. Yarza, P. *et al.* Uniting the classification of cultured and uncultured bacteria and
515 archaea using 16S rRNA gene sequences. *Nature Rev. Microbiol.* **12**, 635-645
516 (2014).
- 517 3. Brown, C. T. *et al.* Unusual biology across a group comprising more than 15% of
518 domain Bacteria. *Nature* **523**, 208-211 (2015).
- 519 4. Castelle, C. J. *et al.* Extraordinary phylogenetic diversity and metabolic versatility
520 in aquifer sediment. *Nat. Commun.* **4**, 2120 (2013).
- 521 5. Iverson, V. *et al.* Untangling genomes from metagenomes: Revealing an
522 uncultured class of marine euryarchaeota. *Science* **335**, 587-590 (2012).
- 523 6. Sekiguchi, Y. *et al.* First genomic insights into members of a candidate bacterial
524 phylum responsible for wastewater bulking. *PeerJ* **3**, e740 (2015).

- 525 7. Wrighton, K. C. *et al.* Fermentation, hydrogen, and sulfur metabolism in multiple
526 uncultivated bacterial phyla. *Science* **337**, 1661-1665 (2012).
- 527 8. Kantor, R. S. *et al.* Small genomes and sparse metabolisms of sediment-
528 associated bacteria from four candidate phyla. *mBio* **4**, e00708-00713 (2013).
- 529 9. Rinke, C. *et al.* Insights into the phylogeny and coding potential of microbial dark
530 matter. *Nature* **499**, 431-437 (2013).
- 531 10. Castelle, C. J. *et al.* Genomic expansion of domain archaea highlights roles for
532 organisms from new phyla in anaerobic carbon cycling. *Curr. Biol.* **25**, 690-701
533 (2015).
- 534 11. Spang, A. *et al.* Complex archaea that bridge the gap between prokaryotes and
535 eukaryotes. *Nature* **521**, 173-179 (2015).
- 536 12. Lynch, M. D. J. & Neufeld, J. D. Ecology and exploration of the rare biosphere.
537 *Nature Rev. Microbiol.* **13**, 217-229 (2015).
- 538 13. Woyke, T. & Rubin, E. M. Searching for new branches on the tree of life. *Science*
539 **346**, 698-699 (2014).
- 540 14. Markowitz, V. M. *et al.* IMG/M 4 version of the integrated metagenome
541 comparative analysis system. *Nucleic Acids Res.* **42**, D568-D573 (2014).
- 542 15. Pruesse, E., Peplies, J. & Glöckner, F. O. SINA: Accurate high-throughput
543 multiple sequence alignment of ribosomal RNA genes. *Bioinformatics* **28**, 1823-
544 1829 (2012).
- 545 16. Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data
546 processing and web-based tools. *Nucleic Acids Res.* **41**, D590-D596 (2013).

- 547 17. Gupta, R. S. The phylogeny and signature sequences characteristics of
548 Fibrobacteres, Chlorobi, and Bacteroidetes. *Crit. Rev. Microbiol.* **30**, 123-143
549 (2004).
- 550 18. Hugenholtz, P., Pitulle, C., Hershberger, K. L. & Pace, N. R. Novel division level
551 bacterial diversity in a Yellowstone hot spring. *J. Bacteriol.* **180**, 366-376 (1998).
- 552 19. Sharp, C. E. *et al.* Humboldt's spa: microbial diversity is controlled by
553 temperature in geothermal environments. *ISME J.* **8**, 1166-1174 (2014).
- 554 20. Costa, K. *et al.* Microbiology and geochemistry of great boiling and mud hot
555 springs in the United States Great Basin. *Extremophiles* **13**, 447-459 (2009).
- 556 21. Cole, J. K. *et al.* Sediment microbial communities in Great Boiling Spring are
557 controlled by temperature and distinct from water communities. *ISME J.* **7**, 718-
558 729 (2013).
- 559 22. Hou, W. *et al.* A comprehensive census of microbial diversity in hot springs of
560 Tengchong, Yunnan Province China using 16S rRNA gene pyrosequencing.
561 *PLoS One* **8**, e53350 (2013).
- 562 23. Varghese, N. J. *et al.* Microbial species delineation using whole genome
563 sequences. *Nucleic Acids Res.* doi:10.1093/nar/gkv657 (2015).
- 564 24. Iino, T. *et al.* *Ignavibacterium album* gen. nov., sp. nov., a moderately
565 thermophilic anaerobic bacterium isolated from microbial mats at a terrestrial hot
566 spring and proposal of *Ignavibacteria* classis nov., for a novel lineage at the
567 periphery of green sulfur bacteria. *Int. J. Syst. Evol. Microbiol.* **60**, 1376-1382
568 (2010).

- 569 25. Podosokorskaya, O. A. *et al.* Characterization of *Melioribacter roseus* gen. nov.,
570 sp. nov., a novel facultatively anaerobic thermophilic cellulolytic bacterium from
571 the class *Ignavibacteria*, and a proposal of a novel bacterial phylum
572 *Ignavibacteriae*. *Environ. Microbiol.* **15**, 1759-1771 (2013).
- 573 26. van der Oost, J., Westra, E. R., Jackson, R. N. & Wiedenheft, B. Unravelling the
574 structural and mechanistic basis of CRISPR-Cas systems. *Nature Rev. Microbio.*
575 **12**, 479-492 (2014).
- 576 27. Makarova, K. S., Wolf, Y. I. & Koonin, E. V. The basic building blocks and
577 evolution of CRISPR-CAS systems. *Biochem. Soc. Trans.* **41**, 1392-1400 (2013).
- 578 28. Inskeep, W. P. *et al.* The YNP Metagenome Project: Environmental parameters
579 responsible for microbial distribution in the Yellowstone geothermal ecosystem.
580 *Front. Microbiol.* **4**, 67 (2013).
- 581 29. Heine, M. & Chandra, S. C. The linkage between reverse gyrase and
582 hyperthermophiles: A review of their invariable association. *J. Microbiol.* **47**, 229-
583 234 (2009).
- 584 30. Weber, K. A., Achenbach, L. A. & Coates, J. D. Microorganisms pumping iron:
585 anaerobic microbial iron oxidation and reduction. *Nature Rev. Microbiol.* **4**, 752-
586 764 (2006).
- 587 31. Aklujkar, M. *et al.* Anaerobic degradation of aromatic amino acids by the
588 hyperthermophilic archaeon *Ferroglobus placidus*. *Microbiology* **160**, 2694-2709
589 (2014).

- 590 32. Holmes, D. E., Risso, C., Smith, J. A. & Lovley, D. R. Genome-scale analysis of
591 anaerobic benzoate and phenol metabolism in the hyperthermophilic archaeon
592 *Ferroglobus placidus*. *ISME J.* **6**, 146-157 (2012).
- 593 33. Hernández-Montes, G., Díaz-Mejía, J. J., Pérez-Rueda, E. & Segovia, L. The
594 hidden universal distribution of amino acid biosynthetic networks: a genomic
595 perspective on their origins and evolution. *Genome Biol.* **9**, R95-R95 (2008).
- 596 34. McCutcheon, J. P. & Moran, N. A. Extreme genome reduction in symbiotic
597 bacteria. *Nature Rev. Microbiol.* **10**, 13-26 (2012).
- 598 35. Giovannoni, S. J., Thrash, J. C. & Temperton, B. Implications of streamlining
599 theory for microbial ecology. *ISME J.* **8**, 1553-1565 (2014).
- 600 36. Cao, B. *et al.* Structure of the nonameric bacterial amyloid secretion channel.
601 *Proc. Natl. Acad. Sci. U.S.A.* **111**, E5439-E5444 (2014).
- 602 37. Evans, M. L. & Chapman, M. R. Curli biogenesis: order out of disorder. *Biochim.*
603 *Biophys. Acta* **1843**, 1551-1558 (2014).
- 604 38. Garcia-Pardo, J. *et al.* Amyloid formation by human carboxypeptidase D
605 transthyretin-like domain under physiological conditions. *J. Biol. Chem.* **289**,
606 33783-33796 (2014).
- 607 39. Hedlund, B. P. *et al.* Potential role of *Thermus thermophilus* and *T. oshimai* in
608 high rates of nitrous oxide (N₂O) production in ~80°C hot springs in the U.S.
609 Great Basin. *Geobiology* **9**, 471-480 (2011).
- 610 40. Murugapiran, S. K. *et al.* *Thermus oshimai* JL-2 and *T. thermophilus* JL-18
611 genome analysis illuminates pathways for carbon, nitrogen, and sulfur cycling.
612 *Stand. Genomic Sci.* **7**, 449-468 (2013).

- 613 41. Giovannoni, S. J. *et al.* Genome streamlining in a cosmopolitan oceanic
614 bacterium. *Science* **309**, 1242-1245 (2005).
- 615 42. Dupont, C. L. *et al.* Genomic insights to SAR86, an abundant and uncultivated
616 marine bacterial lineage. *ISME J.* **6**, 1186-1199 (2012).
- 617 43. Garcia Costas, A. M. *et al.* Complete genome of *Candidatus*
618 Chloracidobacterium thermophilum, a chlorophyll-based photoheterotroph
619 belonging to the phylum Acidobacteria. *Environ. Microbiol.* **14**, 177-190 (2012).
- 620 44. Barns, S. M., Fundyga, R. E., Jeffries, M. W. & Pace, N. R. Remarkable archaeal
621 diversity detected in a Yellowstone National Park hot spring environment. *Proc.*
622 *Natl. Acad. Sci. U.S.A.* **91**, 1609-1613 (1994).
- 623 45. Dick, G. J. *et al.* Community-wide analysis of microbial genome sequence
624 signatures. *Genome Biol.* **10**, R85-R85 (2009).
- 625 46. Nurk, S. *et al.* Assembling single-cell genomes and mini-metagenomes from
626 chimeric MDA products. *J. Comp. Biol.* **20**, 714-737 (2013).
- 627 47. Tennessen, K. *et al.* ProDeGe: a computational protocol for fully automated
628 decontamination of genomes. *ISME J.*, doi:10.1038/ismej.2015.100 (2015).
- 629 48. Stamatakis, A. RAxML-VI-HPC: maximum likelihood-based phylogenetic
630 analyses with thousands of taxa and mixed models. *Bioinformatics* **22**, 2688-
631 2690 (2006).
- 632 49. Ronquist, F. & Huelsenbeck, J. P. MrBayes 3: Bayesian phylogenetic inference
633 under mixed models. *Bioinformatics* **19**, 1572-1574 (2003).
- 634 50. Ludwig, W. *et al.* ARB: a software environment for sequence data. *Nucleic Acids*
635 *Res.* **32**, 1363-1371 (2004).

- 636 51. Pernthaler, A., Pernthaler, J. & Amann, R. Fluorescence in situ hybridization and
637 catalyzed reporter deposition for the identification of marine bacteria. *Appl.*
638 *Environ. Microbiol.* **68**, 3094-3101 (2002).
- 639 52. Yilmaz, L. S., Parnerkar, S. & Noguera, D. R. mathFISH, a web tool that uses
640 thermodynamics-based mathematical models for in silico evaluation of
641 oligonucleotide probes for fluorescence in situ hybridization. *Appl. Environ.*
642 *Microbiol.* **77**, 1118-1122 (2011).
- 643 53. Yilmaz, L. S., Loy, A., Wright, E. S., Wagner, M. & Noguera, D. R. Modeling
644 formamide denaturation of probe-target hybrids for improved microarray probe
645 design in microbial diagnostics. *PLoS One* **7**, e43862 (2012).
- 646 54. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive
647 sequence similarity searching. *Nucleic Acids Res.* **39**, W29-W37 (2011).
- 648 55. Darling, A. E. *et al.* PhyloSift: phylogenetic analysis of genomes and
649 metagenomes. *PeerJ* **2**, e243 (2014).
- 650 56. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for
651 automated alignment trimming in large-scale phylogenetic analyses.
652 *Bioinformatics* **25**, 1972-1973 (2009).
- 653 57. Katoh, K. & Standley, D. M. MAFFT Multiple sequence alignment software
654 version 7: Improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772-
655 780 (2013).
- 656 58. Abascal, F., Zardoya, R. & Posada, D. ProtTest: selection of best-fit models of
657 protein evolution. *Bioinformatics* **21**, 2104-2105 (2005).

59. Huson, D. H. *et al.* Dendroscope: An interactive viewer for large phylogenetic trees. *BMC Bioinform.* **8**, 460-460 (2007).

60. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**, 2460-2461 (2010).

61. Kent, W. J. BLAT—The BLAST-Like Alignment Tool. *Genome Res.* **12**, 656-664 (2002).

62. Brownrigg, R., Minka, T. P., Becker, R. A. & Wilks, A. R. maps: Draw Geographical Maps. R package version 2.1-5. <http://CRAN.R-project.org/package=maps> (2010).

63. Haft, D. H., Selengut, J., Mongodin, E. F. & Nelson, K. E. A guild of 45 CRISPR-associated (Cas) protein families and multiple CRISPR/Cas subtypes exist in prokaryotic genomes. *PLoS Comp. Biol.* **1**, e60 (2005).

64. Zhang, Q., Doak, T. G. & Ye, Y. Expanding the catalog of cas genes with metagenomes. *Nucleic Acids Res.* **42**, 2448-2459 (2014).

65. Makarova, K. S. *et al.* Evolution and classification of the CRISPR-Cas systems. *Nature Rev. Microbiol.* **9**, 467-477 (2011).

66. Bland, C. *et al.* CRISPR recognition tool (CRT): a tool for automatic detection of clustered regularly interspaced palindromic repeats. *BMC Bioinform.* **8**, 209 (2007).

67. Skennerton, C. T., Imelfort, M. & Tyson, G. W. Crass: identification and reconstruction of CRISPR from unassembled metagenomic data. *Nucleic Acids Res.* **41**, e105 (2013).

- 680 68. Alkhnbashi, O. S. *et al.* CRISPRstrand: predicting repeat orientations to
681 determine the crRNA-encoding strand at CRISPR loci. *Bioinformatics* **30**, i489-
682 496 (2014).
- 683 69. Lange, S. J., Alkhnbashi, O. S., Rose, D., Will, S. & Backofen, R. CRISPRmap:
684 an automated classification of repeat conservation in prokaryotic adaptive
685 immune systems. *Nucleic Acids Res.* **41**, 8034-8044 (2013).
- 686 70. Mavromatis, K. *et al.* The DOE-JGI standard operating procedure for the
687 annotations of microbial genomes. *Stand. Genomic Sci.* **1**, 63-67 (2009).
- 688 71. Bushnell, B. BBmap software package, <http://sourceforge.net/projects/bbmap/>
689 (2015).
- 690 72. Margulies, M. *et al.* Genome sequencing in microfabricated high-density picolitre
691 reactors. *Nature* **437**, 376-380 (2005).
- 692 73. Nawrocki, E. P. & Eddy, S. R. Infernal 1.1: 100-fold faster RNA homology
693 searches. *Bioinformatics* **29**, 2933-2935 (2013).
- 694 74. Burge, S. W. *et al.* Rfam 11.0: 10 years of RNA families. *Nucleic Acids Res.* **41**,
695 D226-232 (2013).

698 **Acknowledgments**

699 We thank the DOE JGI production sequencing, IMG, and Genomes OnLine Database
700 teams for their support, along with Steven Quake for metagenomic sequencing and
701 assembly of the Jinze and Gongxiaoshe samples. We thank BC Parks and the Ktunaxa
702 Nation for their cooperation on the Dewar Creek spring. This work was conducted by

the U.S. Department of Energy Joint Genome Institute, a DOE Office of Science User Facility, under Contract No. DE-AC02-05CH11231 and used resources of the National Energy Research Scientific Computing Center, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231. This work was also supported by U.S. Department of Energy (DOE) grant DE-EE-0000716; the U.S. Department of Energy Joint Genome Institute (CSP-182); NASA Exobiology grant EXO-NNX11AR78G; U.S. National Science Foundation grant OISE 0968421; Key Project of International Cooperation by the Chinese Ministry of Science and Technology (MOST, 2013DFA31980); B.P.H. acknowledges generous support from Greg Fullmer through the UNLV Foundation. Metagenome analysis of Dewar Creek was supported in part by funding from Genome Canada, Genome Alberta, Genome BC, and the Government of Alberta (GC Grant 1203). Work at LLNL was conducted under the auspices of DOE Contract DE-AC52-07NA27344 and supported by a Lawrence Fellowship to A. Dekas.

Author Contributions

E.A.E-F., N.C.K., and N.N.I. designed the project; P.F.D., B.P.H., S.E.G., A.L.B., H.D., B.R.B., and W-J.L. provided the samples; J.J., D.G., R.M., and T.W. performed the single-cell experiments; A.E.D. and J.P-R. performed the CARD-FISH experiments; E.A.E-F., D.P-E., J.J., P.F.D., A.P., T.W., N.C.K., and N.N.I. analyzed the data; E.A.E-F., D.P-E., N.C.K, and N.N.I. wrote the manuscript with significant input from P.F.D., B.P.H., T.W., and E.M.R. All authors discussed the results and commented on the manuscript.

726

727 **Competing interests**

728 The authors declare no competing interests.

729

730

731

732 **Figure 1. New lineage identified using metagenomic and single-cell genomic**
733 **approaches.** Workflow used to (A) identify novel SSU rRNA gene sequences globally,
734 along with (B) single-cell genomics pipeline to screen and sequence single cells
735 isolated from geothermal springs samples. For the three geothermal spring
736 environments, we sequenced 13, 2, and 3 SAGs, respectively.

737

738

739

740

741

742

743 **Figure 2. Maximum likelihood concatenated protein phylogeny and cell imaging**
744 **for *Ca. Kryptonia*.** (A) Phylogeny was based on concatenation of 56 conserved marker
745 proteins, where at least 10 marker proteins were used to infer SAG phylogenetic
746 placement (with the exception of JGI-22 with only six marker proteins recovered).
747 Bootstrap support values $\geq 50\%$ are shown with small circles on nodes with robust

phylogenetic support. The *Fibrobacteres-Chlorobi-Bacteroidetes* (FCB) superphylum is shown in the gray shaded region. Expanded phylogenetic tree for *Ca. Kryptonia* shows the placement of the proposed four genera represented by GFMs and SAGs, along with the estimated genome completeness shown in parentheses. **(B)** A *Ca. Kryptonia*-specific FISH (fluorescence *in situ* hybridization) probe was designed and used to visualize cells from Dewar Creek Spring sediment samples. *Ca. Kryptonia* cells hybridizing with the probe are green, while other cells are visualized with 4',6-diamidino-2-phenylindole (DAPI; blue). Scale bar, 5 μ m.

Figure 3. Limited, yet widely dispersed biogeographic distribution of *Ca.*

***Kryptonia* genomes and CRISPR spacers.** All genomic content from the *Ca. Kryptonia* GFMs and SAGs was used to comprehensively search the collection of 640 Gb of assembled metagenomic data from 4,290 environmental samples, including 169 samples from geothermal springs and hydrothermal vents denoted by red triangles (temperature $\geq 50^{\circ}\text{C}$). Marked circles are as follows: (A) Great Boiling Spring, Nevada^{20,21}, (B) Dewar Creek Spring, Canada¹⁹, (C) Jinze pool, Yunnan Province, China²², and (D) Gongxiaoshe pool, Yunnan Province, China²². Significant matches were determined for sequences ≥ 250 bp in length and with $\geq 75\%$ identity threshold for non-ribosomal genomic regions. For metagenomic contigs mapping to the *Ca. Kryptonia* ribosomal operon, a 97% identity threshold was used to capture only high-quality matches to *Ca. Kryptonia*. For CRISPR spacers, only significant matches allowing for up to 3 bp mismatch along the entire length of the spacer were considered. The *Ca.*

Kryptonite genomic hits can be found in Supplementary Table 8 and the manually curated spacer hits can be found in Supplementary Table 7.

Figure 4. Reconstructed metabolic capacity of *Ca. Kryptonite*. Key metabolic predictions and novel features identified in *Ca. Kryptonite* GFM and SAGs, with full gene information available in Supplementary Table 10.

Figure 5. Co-occurrence patterns and metabolic complementarity with *Ca. Kryptonite*. (A) Spearman-rank correlation values were calculated based on reconstructed SSU rRNA sequences across 22 geothermal spring metagenomes, and led to the identification of a cluster of highly correlated phylotypes with *Ca. Kryptonite*. *Armatimonadetes* (cluster 3107) had the highest correlation value ($\rho = 0.82$) with *Ca. Kryptonite*. (B) Biosynthetic pathways present in the *Armatimonadetes* genome which complement missing components in *Ca. Kryptonite*. Each arrow represents an enzymatic component of the biosynthetic pathways; arrows highlighted in blue are contributed by the *Armatimonadetes*, while arrows highlighted in dark orange are contributed by *Ca. Kryptonite*. Black arrows indicate enzyme was not recovered in either.